



Central Analysis Facility



Mark Neubauer

Massachusetts Institute of Technology
for the CDF CAF Development Team

- **System Overview**
- **Milestones/Performance**
- **How do I get my physics done with it?**
- **Future Plans**
- **Conclusions**



The CDF CAF Group



MIT: T.Kim, M. Neubauer, F. Wurthwein

FNAL CD: R. Columbo, G. Cooper, R. Harris, R. Jetton, A. Kreymer, I. Mandrichenko, L. Weems

INFN Italy: S. Belforte, M. Casarsa, S. Giagu, O. Pinazza, F. Semaria, I. Sfligoi, A. Sidoti

Pittsburgh: J. Boudreau, Y. Gotra

Rutgers: F. Ratnikov

Carnegie Mellon: M. Paulini

Rochester: K. McFarland

Special thanks to:

- C.Paus for organizing data server population
- All ProtoCAF/Stage1 commissioning users for much needed feedback!



CAF Computing Model



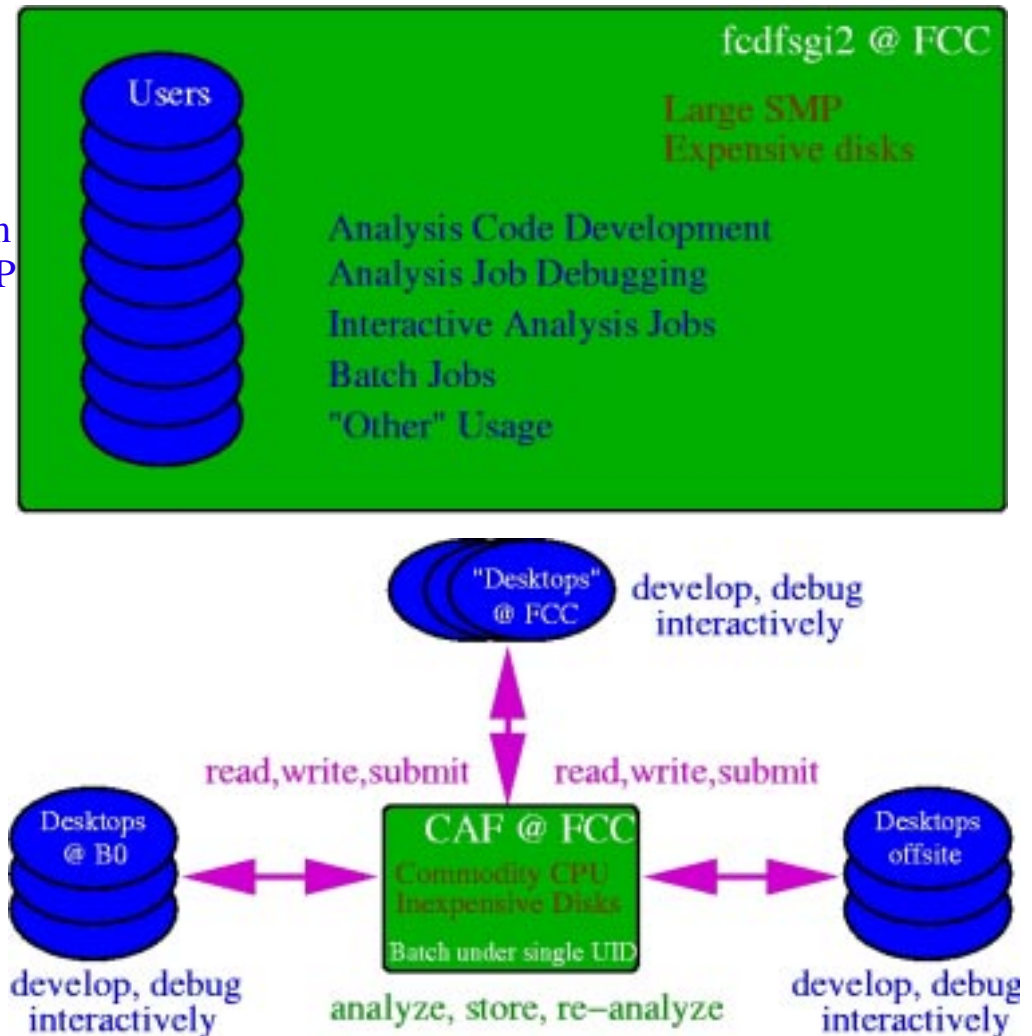
CAF design considerations:

- Submit jobs from 'anywhere'
- Job output can be:
 - sent directly to desktop
 - stored on CAF for later retrieval or input to subsequent job
- Analysis of 5nb dataset in ~1 day for each of 200 users
 - Requires $1\text{THz}/\text{fb}^{-1} \rightarrow$ cheap CPU
- Physics groups want $150\text{TB}/\text{fb}^{-1} \rightarrow$ need cheap disks
 - Replace ~3 IDE drives/week
 - \rightarrow need hot-swap/fault-tol RAID!

See also CDF 5743, 5787, 5802, 5914, 5961

Mark Neubauer/MIT

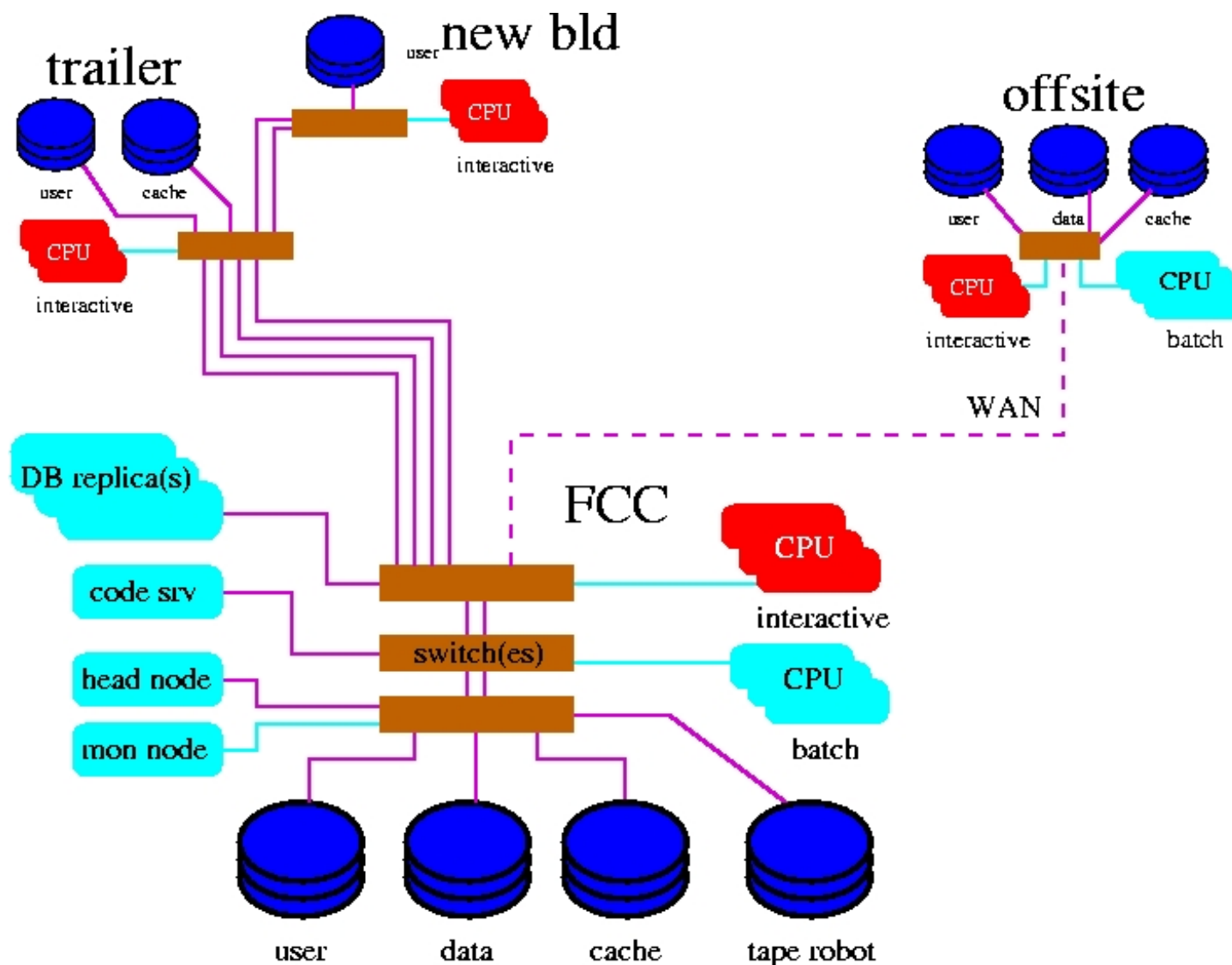
Shift from
large SMP



May'02 CDF Collaboration Meeting



CAF System Implementation





CAF Milestones



- **CDF Central Analysis Computing Review** 8/01-11/01
- **CAF prototype (protoCAF) assembled** 2/25/02
- **Fully-functional prototype system (>99% job success)** 3/6/02
- **ProtoCAF integrated into Stage1 system in FCC** 4/25/02
- **File additional file servers (12TB total) installed** 5/3/02
- **Production Stage1 CAF for collaboration** 5/30/02



ProtoCAF



Stage1





CAF Stage 1 Hardware



Workers (**114 CPUs**, 2U rackmount):

16 Dual Athelon 1.6GHz / 512MB RAM
41 Dual P3 1.26GHz / 2GB RAM
FE (11 MB/s) / 80GB job scratch each

Servers (**35TB total**, 16 4U rackmount):

2.2TB useable IDE RAID50 hot-swap
Dual P3 1.4GHz / 2GB RAM
1 SysKonnnect 9843 Gigabt Ethernet card

Server/Client Performance:

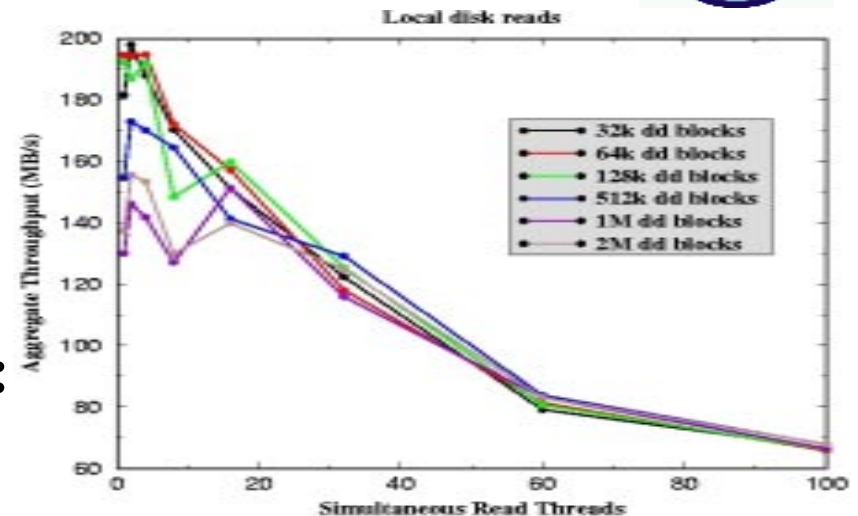
Up to **200MB/s local reads**
~66MB/s NFS for 34 clients (server CPU limited)

For server evaluation details see:

http://mit.fnal.gov/~msn/cdf/caf/server_evaluation.html

or **CDF 5962** (in preparation)

Mark Neubauer/MIT



May'02 CDF Collaboration Meeting



Using the CAF



- Compile, build, debug exe+tcl on 'desktop'
- Adapt CAF shell script example(s) for your job

- Enter appropriate fields & submit job via CAF GUI

CLUI also available

The screenshot shows the 'CDF RunII CAF GUI' window. It has several input fields and buttons. Annotations include:

- A blue arrow pointing to the 'Initial Command' field containing `./simple.sh`.
- A blue arrow pointing to the 'section integer range' field containing `600` and `610`.
- A red oval around the 'Original Directory' field containing `/home/msn/releases/development/CafUtil/examples`.
- A green oval around the 'Output File Location' field containing `msn@fcdlnx2.fnal.gov/cdf/scratch/msn/temp$ tgz`.
- A red arrow pointing to the 'Email Address' field containing `msn@fnal.gov`.
- Buttons for 'Submit', 'Quit', and 'Browse...'.
- A 'Ready' status indicator.
- A log window at the bottom showing job progress:

```
(2002-05-23 01:46:51) Email sent to msn@fnal.gov upon job completion
(2002-05-23 01:46:55) /bin/tar -cvzf /home/msn/msn49959.tgz *
(2002-05-23 01:46:57) Remove /home/msn/msn49959.tgz
(2002-05-23 01:46:57) Job Submission is successful, JID: 873
```

- Monitor job progress or just wait for email from CAF
- Retrieve output using ICAF
- ... or write output directly to 'desktop'!

Web Monitoring of User Queues

Each user a different **queue**

Process type for job length

test: 5 mins

short: 2 hrs

medium: 6 hrs

long: 2 days

This example:

1 job → 11 sections

(+ 1 additional section automatic for job cleanup)

Name	Status	Default Process Type	Share	Prio	Waiting	Ready	Running	Total
akorn	OK	short	1.00	0	0	0	0	0
amitl	OK	short	1.00	0	0	0	0	0
anikeev	OK	short	1.00	0	0	0	0	0
belforte	OK	short	1.00	0	0	0	0	0
msmartin	OK	short	1.00	0	0	0	0	0
msn	OK	short	1.00	0	1	0	11	12
pauly	OK	short	1.00	0	0	0	0	0
paus	OK	short	1.00	0	0	0	0	0
ratnikov	OK	short	1.00	0	0	0	0	0
rescigno	OK	short	1.00	0	0	0	0	0
semeria	OK	short	1.00	0	0	0	0	0
sfiligoi	OK	short	1.00	0	0	0	0	0
sgromoll	OK	short	1.00	0	0	0	0	0
shepard	OK	short	1.00	0	0	0	0	0
sidoti	OK	short	1.00	0	0	0	0	0
spezziga	OK	short	1.00	0	0	0	0	0
test	OK	short	1.00	0	0	0	0	0
thkim	OK	short	1.00	0	0	0	0	0
thom	OK	short	1.00	0	1	0	1	2

Monitoring jobs in your queue

The screenshot shows a Netscape browser window titled "Netscape: FBSWWW - queue msn@CAF". The main content area displays "FBSNG on the web" with a status report for "Farm: CAF", "Time: Thu May 23 01:47:23 2002", and "Report: Queue msn". Below this, there are tabs for "Queues", "Jobs", "Nodes", and "Process Types". The "Jobs" tab is selected, showing a table of jobs with columns: SectID, User, ProcType, Status, Prio, NProc, and Date/Time. The table lists 11 running jobs and 1 pending job. A "User Monitor" section is also visible on the left side of the main window. A blue arrow points from a smaller window on the left to the main window on the right.

User Monitor

Queue Parameters [show]

Status: **OK** Running: 11 Pending: 0

SectID	User	ProcType	Status	Prio	NProc	Date/Time
873.msn_600	cdfcac	short	running	0	1/1	Started at 05/23 01:47:09
873.msn_601	cdfcac	short	running	0	1/1	Started at 05/23 01:47:09
873.msn_602	cdfcac	short	running	0	1/1	Started at 05/23 01:47:10
873.msn_603	cdfcac	short	running	0	1/1	Started at 05/23 01:47:10
873.msn_604	cdfcac	short	running	0	1/1	Started at 05/23 01:47:11
873.msn_605	cdfcac	short	running	0	1/1	Started at 05/23 01:47:11
873.msn_606	cdfcac	short	running	0	1/1	Started at 05/23 01:47:12
873.msn_607	cdfcac	short	running	0	1/1	Started at 05/23 01:47:12
873.msn_608	cdfcac	short	running	0	1/1	Started at 05/23 01:47:12
873.msn_609	cdfcac	short	running	0	1/1	Started at 05/23 01:47:13
873.msn_610	cdfcac	short	running	0	1/1	Started at 05/23 01:47:13
873.msn_end	cdfcac	mailer	waiting	0	0/1	Submitted at 05/23 01:46:57

FCS Group | FBSNG

FBSWWW version 0.1

Monitoring sections of your job

FBSNG on the web
Farm: CAF
Time: Thu May 23 01:47:23 2002
Report: Queue msn

Queues Jobs Nodes Process Types

Queue Parameters [show]

Status: OK Running: 11 Pending: 0

SectID User ProcType

873.msn_601	cdcfcaf	short
873.msn_602	cdcfcaf	short
873.msn_603	cdcfcaf	short
873.msn_604	cdcfcaf	short
873.msn_605	cdcfcaf	short
873.msn_606	cdcfcaf	short
873.msn_607	cdcfcaf	short
873.msn_608	cdcfcaf	short
873.msn_609	cdcfcaf	short
873.msn_610	cdcfcaf	short
873.msn_end	cdcfcaf	mailer

FBSNG on the web
Farm: CAF
Time: Thu May 23 01:48:13 2002
Report: Section 873.msn_600 status

Queues Jobs Nodes Process Types

User Monitor

ID: 873.msn_600 User: cdcfcf

Queue: msn Process Type: short

NProc: 1 Status: running

Need: 0 Depends:

Submitted: 05/23 01:46:57 Started: 05/23 01:47:09

CPU time limit: 2h00m

Proc Rsrc: cpu:100 disk:15 Sect Rsrc:

Command: /fbsng/cafllocal/v1.01/CafExe cdcfcf@fcdhead1.fnal.gov/home/cdcfcf/v1.01/submitter/cafln/msn_%s.tgz msn@fcdflnx2.fnal.gov/cdf/scratch/msn/temp600.tgz msn 4h
cdcfcaf@fcdhead1.fnal.gov/home/cdcfcf/v1.01/submitter/fbs/FBS_%s.msn_600.1.log ./simple.sh 600

Other sections: [msn_600](#) [msn_601](#) [msn_602](#) [msn_603](#) [msn_604](#) [msn_605](#) [msn_606](#) [msn_607](#) [msn_608](#) [msn_609](#) [msn_610](#) [msn_end](#)

(running) (running) (running) (running) (running) (running) (running) (running) (running) (running) (running) (waiting)

Processes

Process #	Node	Status	CPU Time	PID	Command
1	fcdcfcaf057	running	0	6931	CafExe cdcfcf@fcdhead1.fnal.gov/home/cdcfcf/v1.01/submitter/cafln/msn_%s.tgz msn@fcdflnx2.fnal.gov/cdf/scratch/msn/temp600.tgz msn 4h
			0	6940	simple.sh 600
			0	7221	sleep 120

FCS Group | FBSNG

FBSWWW version 0.1



CAF User Tools



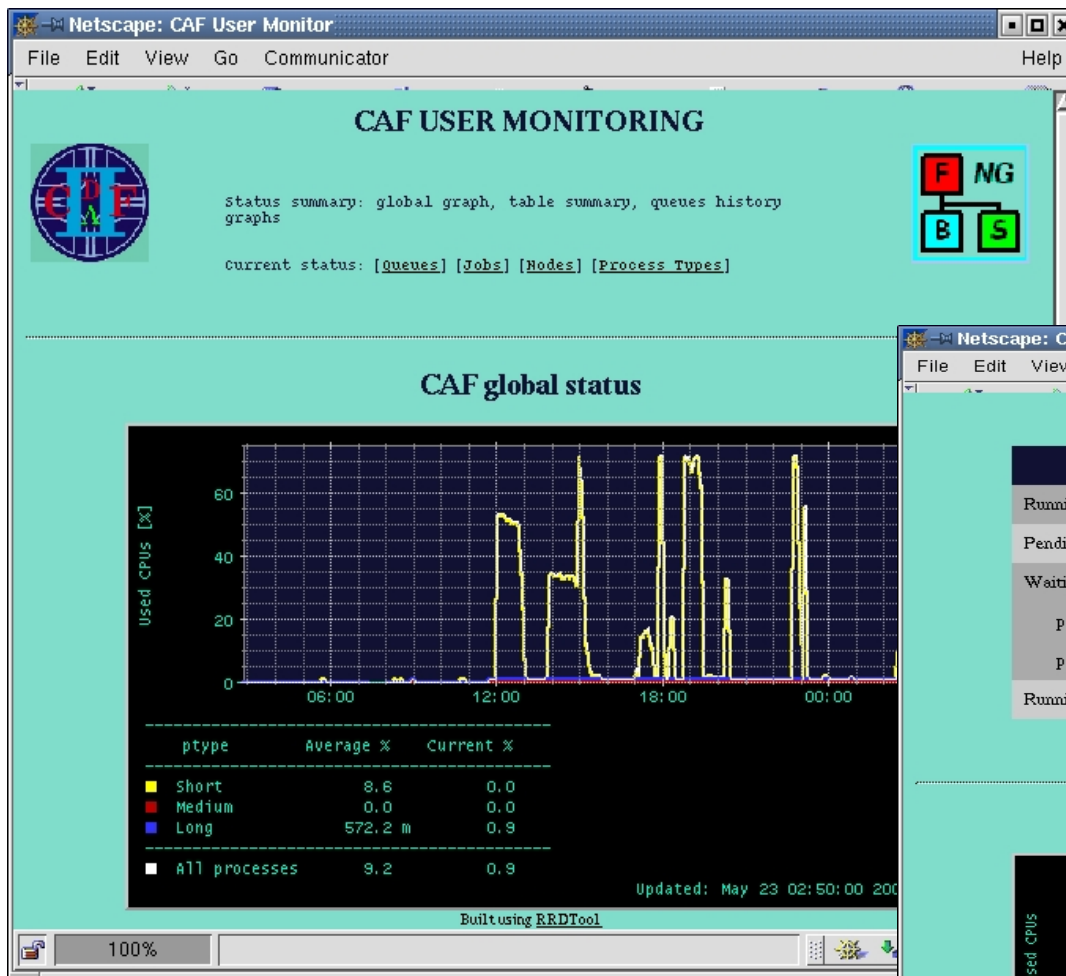
Job Control/Monitoring Utilities:

- **Get list/status of jobs in your queue**
 - > **cafjobs**
- **Check progress of a section**
 - > **caflog** JID sectionNumber
- **Kill a job or sections within a job**
 - > **cafkill** JID [SectionRangeList]

Remote file listings:

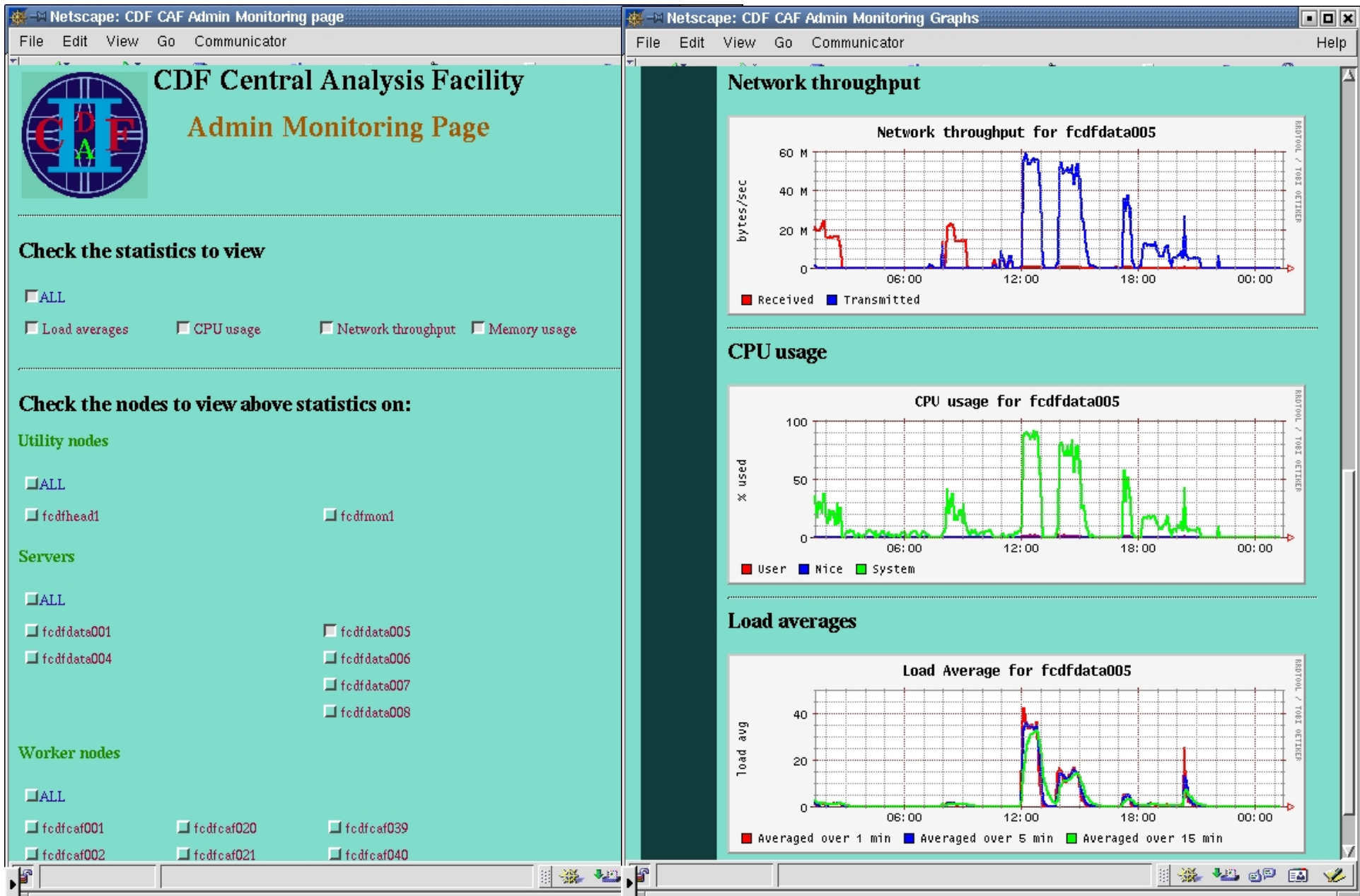
- **Get listing in a section's 'relative path'**
 - > **cafdir** JID sectionNumber [directory]
- **List files in a CAF node's absolute path**
 - > **cafhostdir** hostname directory
- **Get full file listing for a dataset**
 - > **cafdataset** datasetid

Monitoring CAF utilization



...by process type and queue







Job Output Destinations



CAF ftp servers

Job and **scratch** space for each user:

- **Job space** (icaf/) - job output tarballs
Specify *icaf:temp.tgz* for 'Output File Location' in CAF GUI
- **Scratch space** (scratch/) - data for subsequent job(s)
(currently rcp, later kerberized rootd)

**Both spaces accessible via
kerberized ftp and ICAF tools:**

- > icaf_gftp (ICAF GUI)
- > icaf_info (ICAF user-related info)
- > icaf_ls/rm/get

Remote 'desktop'

Output can be sent:

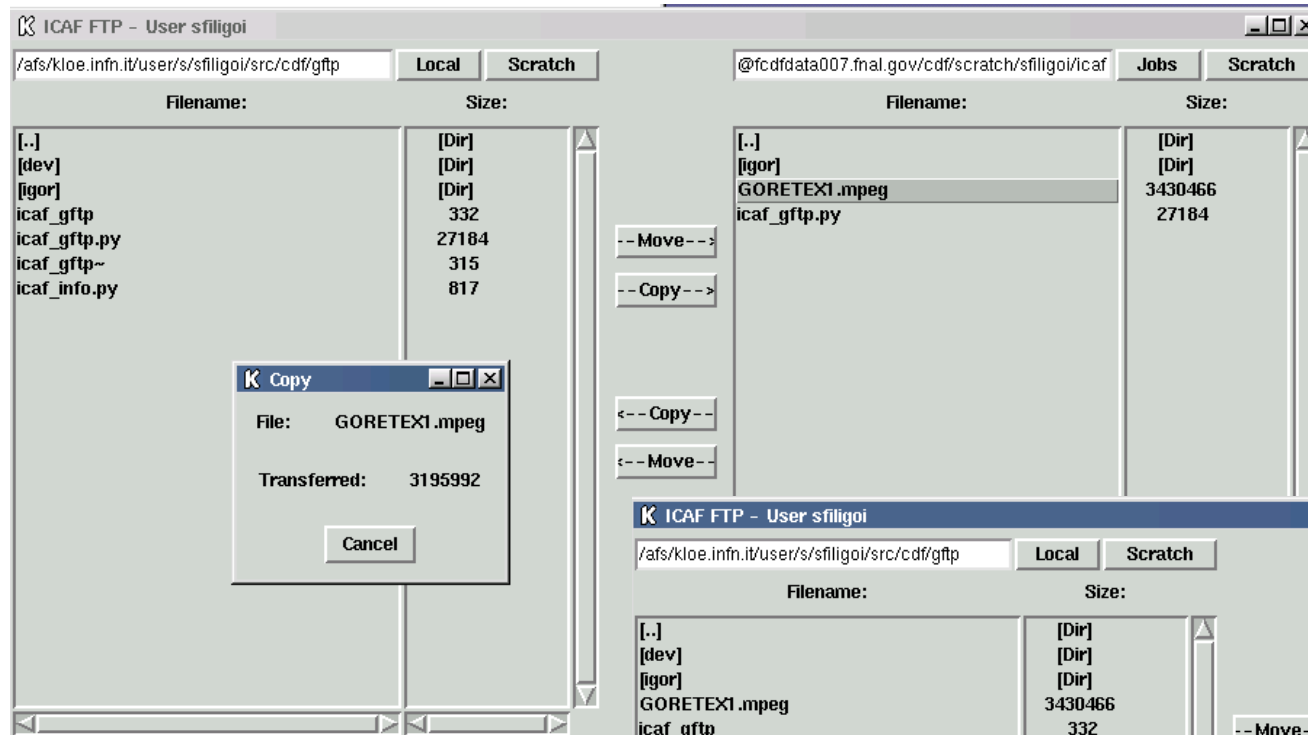
- directly back to desktop
- other remote machine

**Any machine allowing incoming
kerberized rcp can get CAF output**

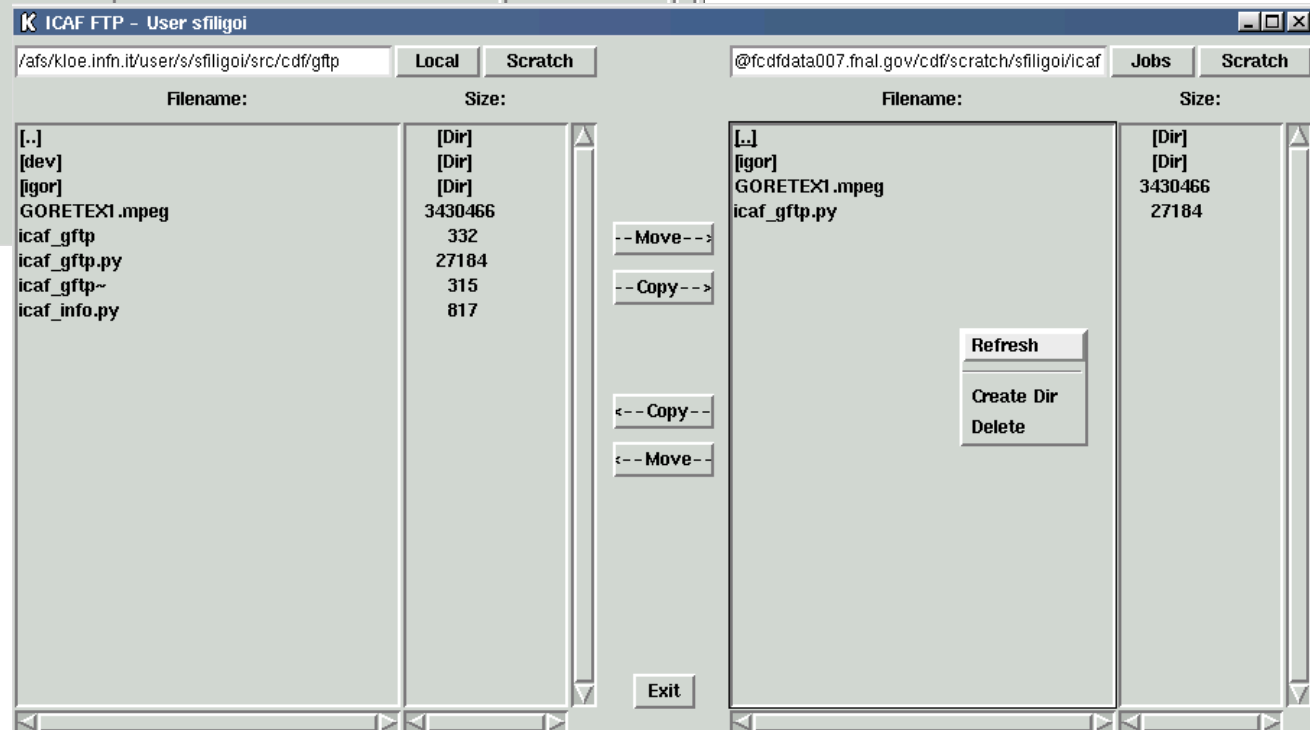
```
> cat ~johndoe/.k5login  
johndoe@FNAL.GOV  
caf/cdf/johndoe@FNAL.GOV
```

special CAF principal for each user

ICAF GUI



remote file retrieval





Data Access on the CAF



Data access methods - NFS & rootd

- **NFS** - server data exported to each worker

talk DHInput

input file /cdf/data/fcdfdata005/jbot0c/*

exit

- **rootd** - remote access to root files

talk DHInput

input file root://fcdfdata005.fnal.gov//export/data/jbot0c/*

exit

Available in DHInput after 4.4.0int6 (no extra work for 4.5.0)

CAF data servers and ftp servers run rootd

→ can run CAF job(s) on output of previous CAF job(s)



Data Management



- **No explicit Data Handling functionality for initial Stage1**
 - Secondary datasets disk-resident on CAF servers (fcdldata001, fcdldata004-006)
- **Datasets on fcdfsi2 mirrored onto CAF servers**
 - Space managed by Strippers Club (C. Paus, et al)
- **Servers currently X% populated with datasets X,Y,Z**
- **Get server data listing using cafdataset/cafhostdir utilities:**
 - > `cafdataset jbot0c`
 - > `cafhostdir fcdldata005 /export/data/jbot0c`

Interfacing to DH big part of evolving Stage 1 CAF!

- Initial Stage1: ~**20 TB** for static data, **9 TB** for dCache development
- Plan: all disk eventually used for dCache/Enstore



Future Plans



Near term:

- **Gain operational experience!!**
 - fix bugs, learn usage patterns & bottlenecks
- **Kerberized rootd**
- **Integrate Data Handling system**
- **Eliminate single points of failure**
- **Develop monitoring/alarms capability** (NGOP, RRD)
- **University ownership issues**
- **Database replicas**

FY03 and beyond:

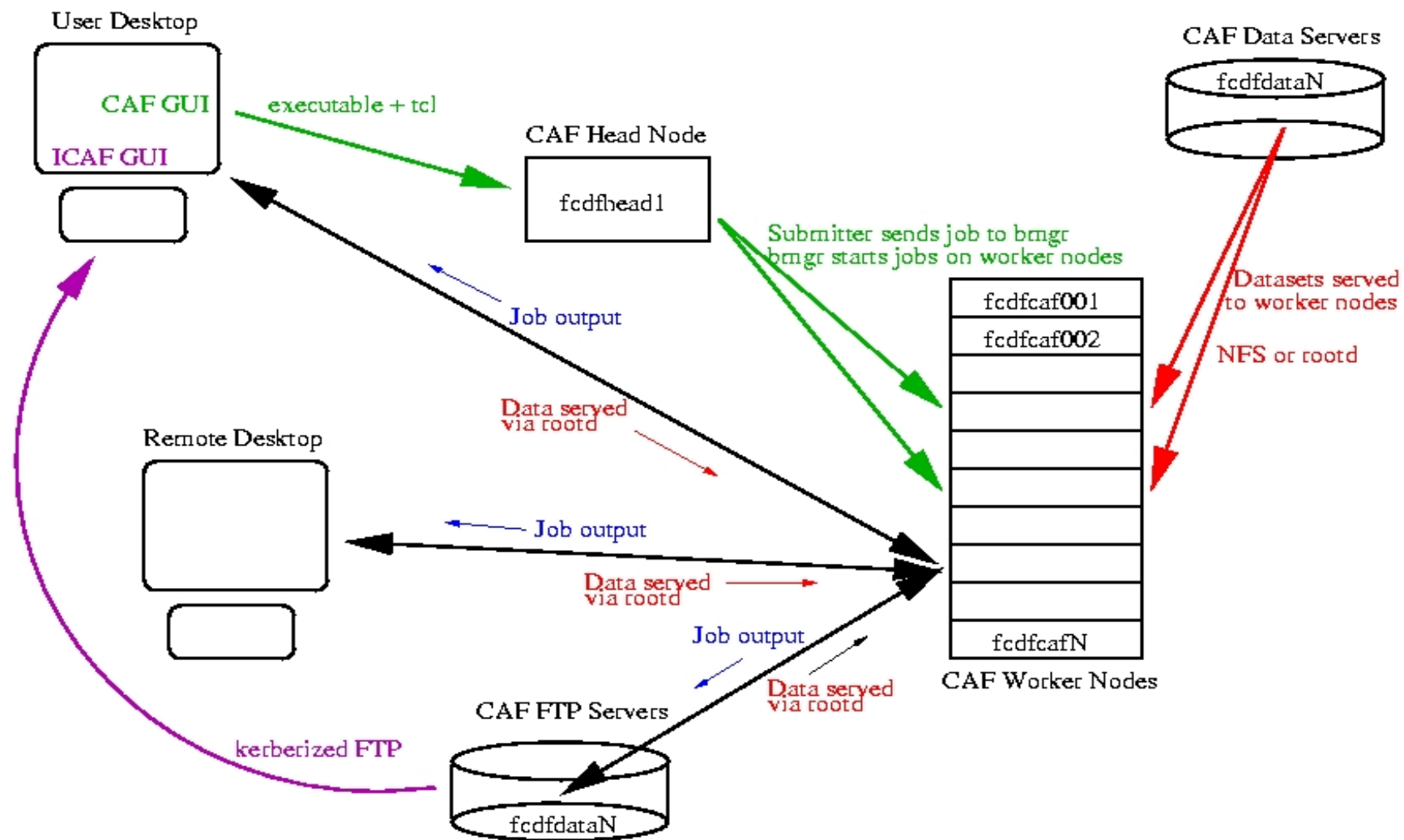
- **Scalability issues**
- **DCAF + GRID**



Summary/Conclusions



- **Paradigm shift in CDF computing**
- **CAF Design→Prototype→Stage1 \leq 6 months**
- **CAF open business!**
 - **Send account request to:**
fkf@fnal.gov + msu@fnal.gov + thkim@fnal.gov
 - **Join CDF_CAF and CDF_CAF_USER mailing lists**
 - **Read the User's Guide**
available at **cdfcaf.fnal.gov** (CAF home page)
 - **Try it out & send CAF-related questions/problems to:**
cdf_caf@fnal.gov





CDF CAF File Server Benchmarking



Remote reads from CAF file server

